

WE CLAIM:

- 1                   1.       A method for identifying a novel nucleic acid molecule encoding a
- 2                   protein of interest comprising:
- 3                    (i)      selecting a specific protein from a first species involved in a
- 4                    regulatory network of interest;
- 5                    (ii)     identifying known proteins that act upstream and
- 6                    downstream in the regulatory network of interest with respect
- 7                    to the specific protein selected;
- 8                    (iii)    constructing the regulatory network of interest from the
- 9                    proteins identified in step (ii);
- 10                  (iv)     for each identified protein, select a domain or motif and
- 11                  search by homology for related proteins in a second species,
- 12                  wherein a related protein is defined as a protein having a
- 13                  homologous domain or motif;
- 14                  (v)      producing a regulatory network for the second species,
- 15                  wherein said regulatory network incorporates the identified
- 16                  related proteins;

17 (vi) comparing the regulatory network from the first species to  
18 the regulatory network of said second species;  
19 (v) identifying a protein present in a regulatory network for one  
20 species but absent in the regulatory network of the other  
21 species; and  
22 (vi) isolating a nucleic acid molecule encoding the protein  
23 identified in step (v) in the species in which it is missing.

1 2. The method of Claim 1 wherein the nucleic acid molecule encodes  
2 human protein.

1 3. The method of claim 1 wherein the related proteins are orthologs.

1  
2 4. The method of claim 1 wherein the regulatory pathway is involved in  
3 apoptosis.

1  
2 5. The method of claim 1 wherein the specific protein from the first  
species is involved in tumor suppression.

1                   6. A method for identifying the affect of a gene knockout on a regulatory  
2 pathway comprising the following steps:

3 (i) identification of the shortest non-oriented pathway  
4 connecting two gene products;  
5 (ii) assigning an initial sign value of “-” to the knockout since the  
6 knockout gene product is inactive;  
7 (iii) moving along the shortest pathway between the two gene  
8 products multiplying the sign with the sign of the next gene  
9 product in the pathway, wherein “-” stands for inhibition, “+”  
10 stands for induction or activation, and “0” stands for the lack  
11 of interaction between two proteins in the specified direction;  
12 and  
13 (iv) determining the final sign at the end of the pathway, wherein  
14 “-” indicates inhibition and “+” indicates induction or  
15 activation of the pathway.

7. A method for identifying a novel nucleic acid molecule encoding a protein of interest comprising:

- (i) selecting a gene of interest and searching a database for homologous sequences;
- (ii) aligning the homologous sequences identified in step (i);
- (iii) constructing a gene tree using the sequence alignment;
- (iv) constructing a species tree;
- (v) imputing the species tree and gene tree into an algorithm which integrates the species tree and the gene tree into a reconciled tree; and
- (vi) identifying orthologous genes present in one species but missing in another.

1                   8. The method of claim 7 wherein the following algorithm is used to  
2 integrate the species tree and the gene tree into a reconciled tree:

3                   (i) computing the similarity  $\sigma(S_{gi}, S_{sj})$  for each pair of interior  
4                   nodes from trees  $T_g$  and  $T_s$ ,

5                   (ii) finding the maximum  $\sigma(S_{gi}, S_{sj})$ ;

6                   (iii) saving  $S_{gi}$  as a new cluster of orthologs, save  $\{S_{gi}\} - \{S_{sj}\}$  as  
7                   a set of species that are likely to have gene of this kind (or  
8                   lost it in evolution);

9 (iv) eliminating  $S_{gi}$  from  $T_g$ ;  $T_g' = T_g \setminus S_{gi}$ ;

10 (v) repeating step (ii)-(iv) until  $T_g$  is non-empty.

11 9. A method for identifying a novel gene comprising the following  
12 steps:

13 (i) defining a motif or domain composition of a gene of interest;

14 (ii) searching for sequences which correspond to nucleotide  
15 sequences in an expression sequence tag database or other  
16 cDNA databases using a program such as BLAST and  
17 retrieving the identified sequences;

18 (iii) searching additional databases for expressed sequence tags  
19 containing the domains and motifs characteristic for  
20 the gene of interest with Hidden Markov Model of domains  
21 and motifs identified in step (i);

22 (iv) identifying nucleotide sequences comprising the gene of  
23 interest.

24 10. The method of claim 9 further comprising using each identified  
25 expression sequence tag to search sequence databases for

26 overlapping sequences for the purpose of assembling longer  
27 overlapping stretches of DNA.

28  
29      *Su* *5* 11. A method for extracting information on interactions between  
30      biological entities from natural-language text data, comprising:  
31              (i) parsing the text data to determine the grammatical structure of the  
32                      text data, and  
33              (ii) regularizing the parsed text data to form structured word terms.

1                   12. The method according to claim 11, further comprising preprocessing  
2                   the data prior to parsing, with preprocessing comprising the step of identifying biological  
1                   entities

1                   13. The method according to claim 11, further comprising referring to an  
2 additional parameter which is indicative of the degree to which subphrase parsing is to be  
1 carried out.

1                           14.     The method according to claim 11, wherein said parsing step further  
2     comprises segmenting the text data by sentences.

1           15. The method according to claim 11, wherein said parsing step further  
2 comprises:

3           segmenting the text data by sentences; and  
4           segmenting each of the sentences at identified words or phrases.

1           16. The method according to claim 11, wherein said parsing step further  
2 comprises:

3           segmenting the text data by sentences; and  
4           segmenting each of the sentences at a prefix.

1           17. The method according to claim 11, wherein said parsing step further  
2 comprises skipping undefined words.

1           18. The method according to claim 11, wherein said parsing step further  
2 comprises:

3           identifying one or more binary actions and their relationships; and  
              identifying one or more arguments associated with the actions.

1      *sub b* }      19.     The method according to claim 11, further comprising performing  
2     error recovery when parsing of the text data is unsuccessful.

1                      20.     The method according to claim 19, wherein said error recovery step  
2     comprises:

3                      segmenting the text data; and  
4                      analyzing the segmented text data to achieve at least a partial parsing of the  
5     unsuccessfully parsed text data.

1      *sub d1* }      21.     The method according to claim 11, wherein said tagging step  
2     comprises providing the structured data component in a Standard Generalized Markup  
1     Language (SGML) compatible format.

1                      22.     A computer system for extracting information on biological entities  
2     from natural-language text data, comprising:  
3                      (i)     means for parsing the natural-language text data; and  
4                      (ii)    means for regularizing the parsed text data to form structured word  
5                      terms.

1           23.     The system according to claim 22, further comprising means for  
2     preprocessing the data prior to parsing, with the preprocessing means comprising  
3     identifying biological entities.

1           24.     The system according to claim 22, further comprising means for  
2     referring to an additional parameter which is indicative of the degree to which subphrase  
1     parsing is to be carried out.

1           25.     The system according to claim 22, wherein said parsing means  
2     further comprises means for segmenting the text data by sentences.

1           26.     The system according to claim 22, wherein said parsing means  
2     further comprises:  
3                 means for segmenting the text data by sentences; and  
4                 means for segmenting each of the sentences at identified words or phrases.

1           27.     The system according to claim 22, wherein said parsing means  
2     further comprises:  
3                 means for segmenting the text data by sentences; and

4 means for segmenting each of the sentences at a prefix.

1 28. The system according to claim 22, wherein said parsing means  
2 further comprises means for skipping undefined words.

1 29. The system according to claim 22, wherein said parsing means  
2 further comprises:  
3 means for identifying one or more binary actions and their relationships; and  
4 means for identifying one or more arguments associated with the actions.

1 30. The system according to claim 22, further comprising means for  
2 performing error recovery when parsing of the text data is unsuccessful.

1 31. The system according to claim 22, wherein said error recovery  
2 means comprises:  
3 means for segmenting the text data; and  
4 means for analyzing the segmented text data to achieve at least a partial  
5 parsing of the unsuccessfully parsed text data.

1                   32.     The system according to claim 22, wherein said tagging means  
2     comprises means for providing the structured data component in a Standard Generalized  
3     Markup Language (SGML) compatible format.

1

CONFIDENTIAL - SECURITY INFORMATION - THIS DOCUMENT CONTAINS INFORMATION WHICH IS UNCLASSIFIED